

An Initial Study on Ideal GUI Test Case Replayability

Arthur-Jozsef Molnar

Department of Computer Science

Faculty of Mathematics and Computer Science

Babeş-Bolyai University

Abstract—In this paper we investigate the effect of long-term GUI changes occurring during application development on the reusability of existing GUI test cases. We conduct an empirical evaluation on two complex, open-source GUI-driven applications for which we generate test cases of various lengths. We then assess the replayability of generated test cases using simulation on newer versions of the target applications and partition them according to the type of repairing change required for their reuse.

I. INTRODUCTION

Many software applications today employ graphical user interfaces (GUIs) to interact with users. As a highly successful paradigm, we encounter GUI-driven applications on many types of devices, a trend that seems set to continue in today's world of pervasive computing. However, while these applications are ubiquitous, the same thing cannot be said about the processes that should support their life cycle, such as quality assurance (QA). Because as much as 50% of application code can be GUI related [1], the existence of QA processes for these applications becomes crucial. Recent work on GUI testing consists of notable contributions such as developing theoretical frameworks for GUI testing [1]–[3], implementing advanced testing tools [4], [5], studying possibilities for automation [6], [7] and gathering empirical evidence regarding the success of automated testing of complex applications [6], [8], [9].

Tools for GUI testing can be divided into two categories: capture-replay and model-based.

Capture-replay tools [10] such as Pounder¹, Marathon² or jfcUnit³ are generally representative of the first wave of automated tooling and work in the two phases that spawned their name: during the first capture phase, the tester works with the application under test (AUT) and manually records test cases which are stored by the tool and then replayed during the second phase. This approach introduces the GUI paradigm to the generation of test cases by allowing them to be built by interacting with the AUT. However, typical capture-replay tools suffer limitations when the behaviour or GUI of the tested systems change. Also, it must be noted that such tools are only able to automate one part of the testing process, as creating test cases remains an overwhelmingly manual undertaking. Also, typical capture-replay tools cannot

provide comprehensive oracles for GUI testing [11] beyond crash-detection and recognizing error windows. These issues were well known to practitioners and tool developers so many ideas were implemented to alleviate such limitations.

In [12], Takahashi proposes intercepting Win32 API graphic calls to replace screen captures, as they are more reliable and occupy less storage space when persisted. Another approach is described in [13], where Ostrand et al. design a Test Development Environment (TDE) that links a test designer and a test generation library with a standard capture-replay tool. Its feasibility is then tested by generating test cases for a medical diagnosis machine.

Many of the shortcomings of capture-replay tools can be addressed using model based approaches. The availability of a model allows automated generation and execution of test cases and helps with implementing necessary test oracles to evaluate testing results. The typical drawback is represented by the time and effort dispensed for model building and validation, as a suitable balance must be achieved between model complexity and system testability to enable revealing system faults [14]. Recent advances addressing such issues come from Microsoft Research's NModel, that uses C# for building the model [15] and Silva et al. [16] who employ Spec Explorer for GUI testing.

A solid body of research in model-based GUI testing⁴ was initiated by Memon's PhD thesis [1]. He provides a definition for the state of the GUI ([1], p.29), the event-flow graph ([1], p.37) which models the valid flow of GUI events and for the GUI test case ([1], p.42). The theoretical foundations are then used for implementing the GUITAR⁵ GUI testing framework used in the presented research.

One of the outstanding issues in automated testing regards test case maintenance. Like all application layers, GUIs invariably change during application development and maintenance. Widgets can be resized, moved or changed in numerous ways to fit new application requirements. This leads to many test cases becoming unusable on newer versions of the target applications as testing tools cannot recognize the changed GUI elements. This constitutes a major hurdle in the automation of the process, regardless of approach.

Our research aims to assess how typical changes in GUI-

¹<http://pounder.sourceforge.net>

²<http://www.marathontesting.com>

³<http://jfcunit.sourceforge.net>

⁴Literature refers to *GUI testing* as testing an application through its GUI.

⁵<http://guitar.cs.umd.edu>

driven applications affect the reusability of existing test cases. In this sense, we study an idealized situation where we categorize existing test cases according to their degree of replayability using known correct information obtained by studying GUI changes occurring in two complex, open-source GUI applications.

The structure of this paper is as follows: the second section introduces required preliminaries and related work. The third section details the target applications used in our research and presents our initial case study which provides both cross-sectional and longitudinal perspectives. The fourth section overviews threats to the validity of our empirical evaluation while the last section is reserved for conclusions and future work planned.

II. PRELIMINARIES

In this section we describe related work regarding GUI test case maintenance and we briefly describe the GUITAR GUI testing framework that we extensively use in our case study.

A. Related work

The problem of GUI test case maintenance has not gone unnoticed and several approaches have been proposed. In [17], Memon proposes *repairing transformations* that insert and remove test case steps with the aim of repairing tests broken by changes in the AUT. An empirical evaluation is then performed where the efficiency of the proposed process is evaluated using four open source applications, one of which is the FreeMind mind-mapper also employed in our research. Huang et al. use genetic algorithms to repair automatically generated infeasible test cases in [18]. In their study they use the same theoretical foundation as [17] and evaluate obtained results on several synthetic applications. These approaches prove that GUI test cases can be successfully repaired to run on modified versions of the AUT. However, they are of limited use in long-term regression testing because by altering the sequence of test steps they do not replay *exactly* the same test steps on the modified application. McMaster and Memon detail preliminary work in enabling regression testing of GUI applications in [19], where they describe a conceptual heuristic process to find functionally equivalent widgets across versions of a GUI application. In our previous research we implemented such a process and showed it achieves high accuracy in correctly classifying GUI elements across many application versions [20].

However, because all previous approaches to repair test cases are subject to error, we were interested in performing an evaluation on the efficiency of an ideal error-free process when employed for long-term GUI test case maintenance in the case of complex GUI applications. Such an evaluation provides a benchmark against which existing and future implementations can be compared and is useful for assessing how many test cases can be replayed when using perfect⁶ implementations. Such a perfect implementation can be obtained by formally

documenting GUI changes, annotating GUI elements so they are recognized across versions by the testing harness or approximated using highly accurate heuristic processes.

B. The GUITAR framework

The theoretical aspects presented in [1] were implemented in the GUITAR testing framework [4], [6]. GUITAR is a mature testing toolset that automates many processes in testing: obtaining the GUI model, generating valid test cases, replaying them on the AUT and recording information usable by test oracles. Therefore it facilitates the adoption of model based testing by decreasing the effort of obtaining the GUI model and providing automation for associated activities. GUITAR's components are available⁷ for Java and Web, with Windows, Android and iOS implementations currently in development. This makes GUITAR a state of the art tool in the research and practice of GUI application testing. Its four main components, briefly described in the order in which they are usually employed are:

- *GUIRipper* can be used to automatically record the AUT's GUI model in XML format. Using reflection and automated interaction, it records all accessible application windows together with their widgets and properties [4].
- *GUI2EFG* takes as input the model obtained using the GUIRipper and computes the application's event-flow-graph that provides the valid event sequences within the GUI of the AUT [21]. This component provides crucial functionality for building valid test cases, as not every widget is actionable at all times.
- *TestCaseGenerator* can build valid test cases using the provided GUI model and event-flow-graph. The *TestCaseGenerator* uses a plugin architecture that allows implementing new strategies for generating test cases [7].
- *TestCaseReplayer* is used to run generated test cases and record the target GUI state after each test step, allowing offline analysis of test case execution. The *Replayer* component was used in several studies. In [6] authors use it for research in regression testing, while in [22] Brooks and Memon employ it for automating profile-guided testing. In [23], Xie and Memon use GUITAR to study desirable characteristics of GUI test suites while in [24] a pilot study assesses GUI event interactions and the influence of event context on test case outcome.

Some of GUITAR's limitations stem from the theoretical foundations from which it was developed [1] and regard testing GUIs that present continuous streams of data or that interface real-time systems. Due to their particular constraints, they might not be handled properly by GUITAR's components and are therefore not targeted by our research. More so, to the best of our knowledge there exist no readily-available tools for assisting QA processes targetted towards such GUIs, leaving GUITAR as the prime candidate for carrying out our empirical investigation. In our research we use *GUIRipper* to obtain the GUI models of our target applications, which we then pass

⁶Repaired test cases are as functionally close to the original as possible.

⁷As of February, 2012

through *GUI2EFG* to compute the valid event sequences. Finally, we use GUITAR’s test generation component to generate the test cases used in our evaluation.

III. CASE STUDY

In this section we present our initial case study in which we investigate how changes that occur during application development affect replayability of existing GUI test cases. As GUI test case steps action widgets, assessing test case replayability requires information about the functionally equivalent widgets [19] between the studied versions. This information was obtained using an automated heuristic process [20], with all results manually double-checked for correctness. The amount of effort involved limited our study to two target applications: the FreeMind⁸ mind-mapper and the jEdit⁹ text editor, both detailed in the following section. Using their publicly available source code repositories we downloaded 30 distinct versions of these applications to obtain a suitable balance between generality and the amount of effort required to prepare the data. Next, we generated comprehensive test suites for each version using GUITAR’s test generator. One of our goals was assessing the effect of GUI changes on test cases of different lengths. Therefore, for each application version we generated a test suite in the following manner:

- *Event-interaction coverage*¹⁰. We generated all such test cases. Running such tests was proposed in previous work by Xie [8] as a method to provide automated crash-testing. The number of obtained test cases varied between target applications. In the case of the FreeMind versions, the number of event-interaction test cases varied between 825 and 9,175. For jEdit, the number of length-2 test cases varied between 3,453 and 9,601.
- *Randomly generated length 3,4 and 5 test cases*. Generating all length- n , $n > 2$ test cases is not feasible for complex applications. For example, the number of length-3 test cases varied between 12,555 for the simplest version of FreeMind and 323,211 for the most complex jEdit version. Longer sequences increase exponentially in count: for jEdit version 4.3.0final over 300 million length-5 test cases can be generated. Our strategy was therefore to generate 10,000 random test cases for each test case length using GUITAR’s *RandomSequenceLengthCoverage* test generator plugin. This allows keeping the number of test cases reasonable while properly sampling the target application’s event flow graph.

Our strategy resulted in 404,826 FreeMind and 564,869 jEdit test cases that we believe properly sample the target applications’ test case space. The next step was to classify generated test cases in one of four categories for all subsequent versions of the target application. The categories were considered to be representative both for state of the art testing tools such as GUITAR and for more advanced implementations that

are able to repair or update test cases, such as ones proposed by Memon [17], Huang [18], or our own proposed approach [20] based on McMaster and Memon’s preliminaries [19].

In the following we detail the four categories:

- 1) *Replayable using widget Id*. This situation simulates how test cases can be replayed using tools that identify widgets using assigned Id’s, an example of which is the GUITAR framework itself. GUITAR components use widget properties to calculate the *Ids* that are reused when replaying test cases. This approach is generally more accurate than the first wave of capture-replay tools, many of which use positional information to find GUI widgets.
- 2) *Replayable after repair*. This category is comprised of test cases that can be repaired using previously described approaches [9], [20], [25] to be replayable on the newer version of the application. This amounts to all GUI elements actioned in the test case, including reaching steps required for enabling further GUI actions to have equivalents on the newer version and their sequence to remain valid according to the newer version’s event flow graph. Of course, all test cases that are replayable using widget *Ids* are also replayable using a hypothetical identity repair that does not perform any changes. Test cases in these first categories are exactly replayable on the newer application version provided that highly accurate processes are implemented.
- 3) *Repairable*. Compared to the category described above, we relax the imposed event flow graph condition and we only require the existence of equivalent widgets on the new application version. Test cases in this category will not have the same sequence of events as the original ones, but they can be repaired using approaches detailed in [17] or [18].
- 4) *Unrepairable*. This last category comprises test cases that cannot be repaired. This is due to at least one of the test case widgets missing from the newer application’s GUI. To the best of our knowledge, the only way of salvaging these test cases is adding or removing test steps, as detailed by Memon [17].

Our goal is to evaluate the replayability of GUI test cases on newer versions of the target applications. For this, we considered the generated test cases for all examined versions and categorized them for all subsequent application versions. For example, the 30,824 test cases generated for FreeMind 0.2.0¹¹ were categorized for all 12 subsequent examined versions of the application, the last of which is a September 2007 CVS snapshot of FreeMind.

The present section contains three subsections. The first one details our chosen target applications. The following two subsections present our cross-sectional and longitudinal evaluations. We performed the cross-sectional study to assess the replayability of test suites on consecutive application versions, while the longitudinal approach overviews the results obtained

⁸http://freemind.sourceforge.net/wiki/index.php/Main_Page

⁹<http://jedit.org>

¹⁰All valid length-2 test cases.

¹¹The first FreeMind version examined

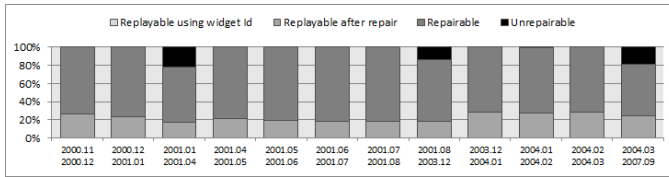


Fig. 1. Cross-sectional FreeMind test case replayability

over all the studied versions. In order to limit the effect of randomness in test case generation the procedure was repeated three times with average values obtained reported.

A. Target applications

The chosen applications for our case study are two complex, open source GUI-driven applications available under non-restrictive GPL-licences: the FreeMind mind mapper and the jEdit text editor. We chose 30 distinct versions of these applications that range between November 2000 and September 2007 for FreeMind¹² and January 2000 to May 2010 for jEdit¹³. Both applications are available on SourceForge¹⁴ where they rank among the most popular applications: FreeMind recorded over 14.3 million downloads over its lifetime while jEdit was downloaded over 6.7 million times since the project was started. Both applications received the "Project of the Month" SourceForge award over their lifetimes. With regards to complexity, FreeMind's GUI consists of one main window having between 101 and 280 GUI elements, while jEdit contains between 12 to 16 windows that contain between 482 and 992 GUI elements. Due to limitations of the GUIRipper tool, the "Options" window of both applications was not recorded and is disregarded in the present study. A detailed overview of the studied application versions is available in [26].

B. Cross-Sectional Approach

The first part of our investigation was undertaken to obtain a detailed picture of GUI test case replayability over targeted application versions. For this, we categorized test cases generated for each version according to their replayability on the application version immediately following it. This approach can isolate application versions for which test cases cannot be repaired and allows corroborating known changes in the GUI with their effect on test case replayability.

Figure 1 shows our obtained results for FreeMind. As the application source code was obtained directly from CVS we use timestamps to identify application versions. The first immediate observation is that the number of directly replayable test cases is of no significance. However, we also observe that most test cases can be repaired to run on the newer version of the application. Most unrepairable test cases are found in April 2001, December 2003 and September 2007 versions,

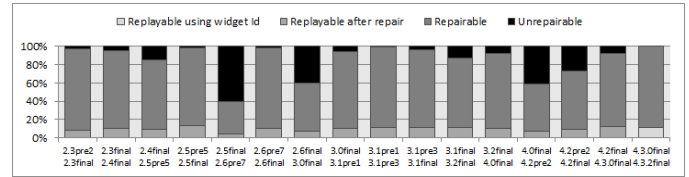


Fig. 2. Cross-sectional jEdit test case replayability

which according to our data [26] correspond with major GUI changes.

The corresponding results for jEdit are shown in Figure 2. Again we witness most generated test cases being at least repairable across version pairs. However, we must note that several version pairs exhibit large numbers of unrepairable tests. The best example is between versions 2.5final and 2.6pre7, where 60% of tests become unrepairable. This can be explained using our jEdit GUI model that consists of 16 windows for the former version and 12 for the latter. Tests that target eliminated windows naturally become unrepairable, as the target section of the AUT no longer exists.

Taking into account the long time span of the examined versions (7 years for FreeMind and 10 for jEdit), together with jEdit's complex GUI [26], we conclude that the major factor affecting test case repairability is represented by functional changes in the application GUI. We believe these initial results confirm Memon's previous findings [17] and hint toward the usefulness of efficient approaches for repairing GUI test cases.

C. Longitudinal Approach

In the second part of our investigation we examine long-term replayability of GUI test cases. For this, we categorize test cases generated for each version on all the subsequent versions of the same application. This provides insight into the long term effects GUI changes have on existing test cases.

Figure 3 shows the results of our longitudinal evaluation on the FreeMind application. Each column group shows replayability in the last studied version of FreeMind for test cases generated in the given version. For example, as the last studied version of FreeMind has a timestamp of September 2007, the second column shows the replayability of test cases generated for the December 2000 version on the latest snapshot. As our information only includes data on functionally equivalent widgets in consecutive versions, we had to categorize test cases on each of the intermediary versions up to the final one. In our example, we had to categorize the test suite generated for the December 2000 version on the 10 intermediate versions separating it from the final one. The four columns in each group symbolize test cases of lengths 2,3,4 and 5 in left to right order.

The data confirms our educated guess that the "age" of test cases has an important effect on their replayability. However, we observe that roughly 50% of all test cases remain at least repairable after 7 years of application development. We attribute this partly to FreeMind's simpler user interface that

¹²13 intermediate versions between 0.2.0 and 0.8.0

¹³17 versions between 2.3pre2 and 4.3.2final

¹⁴<http://sourceforge.net>

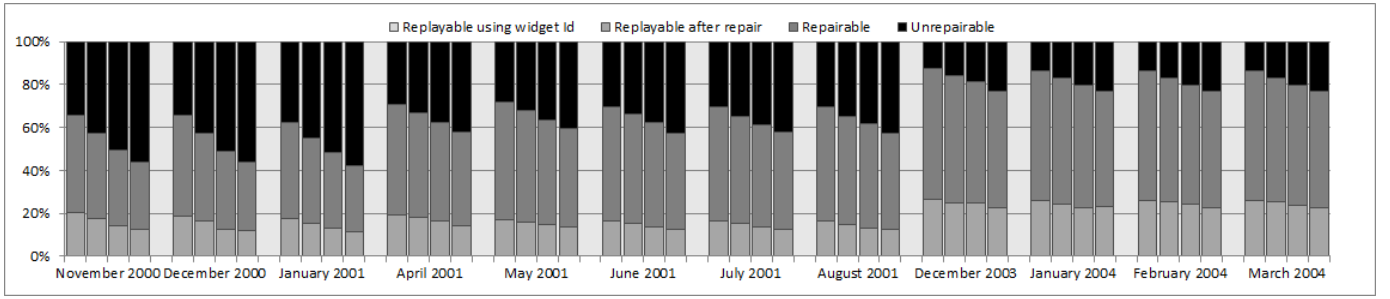


Fig. 3. Longitudinal FreeMind test case replayability

consists of one window¹⁵ that contains all its widgets. Also, similar to Memon [17], we find that test case length has a considerable effect on replayability, as longer test cases are more prone to becoming unrepairable. If only one of a test case's events lacks a functional equivalent on the newer version it is immediately categorized as unrepairable. This leads to an interesting issue that may appear when given three consecutive application GUI versions, say G_1 , G_2 and G_3 . Let us assume that event $e \in G_1$, $e \notin G_2$ but $e \in G_3$, so a GUI event that is removed from an intermediary version but reappears at a later date. This might cause test cases containing e , which are unrepairable on G_2 , to become at least repairable for G_3 . As generally no information is available regarding an application's *future* GUI structure, we consider unrepairable test cases to remain as such. This also holds for the other two categories employed.

Figure 4 provides the longitudinal data for the jEdit application. Our observations for FreeMind generally hold in the case of our second application: older or longer test cases have higher probability of not being reusable. An interesting aspect regards the result of jEdit's GUI model changes between versions 2.5final and 2.6pre7. In the previous section we witness GUI changes between the versions making roughly 60% of test cases unrepairable on the latter version. In the longitudinal view the specified version pair acts as a choke-point, as we witness most test cases generated prior to version 2.6pre7 being unusable on our latest studied version, 4.3.2final, due to them being broken by changes in version 2.6pre7. This enforces our belief that major GUI changes are a serious hurdle in automating GUI regression testing, and processes based on research such as detailed in [17], [18], [20] must take steps to alleviate these issues. Also, our evaluation shows that investing effort in processes related to test case maintenance is worthwhile, as a highly accurate automated process will be able to consistently repair old test cases to work on new versions of the AUT, even with long timespans considered.

IV. THREATS TO VALIDITY

We partition threats to the validity of our empirical evaluation into internal and external. Internal threats are represented by errors in the process employed to obtain our data. In this

regard, the main issue regards the fact that our evaluation was performed by *simulating* test case execution using test case data together with GUI and event flow graph models. This enabled us to generate and categorize over 24 million test cases, including ones generated for intermediate versions for the longitudinal study. However, this has the drawback that errors in recorded GUI models or implementation peculiarities of GUITAR components or the applications themselves might cause test cases to belong to other categories than those assigned to in our study.

External threats regard the generalization of obtained information. While FreeMind and jEdit are complex real-life applications having extensive user-bases and participating in previous empirical research [8], [9], [17], [22], [24], they are not representative of all possible GUI implementations. Practitioners looking to capitalize on our results must have a good understanding on the presented limitations and peculiar aspects regarding their targeted applications.

Our best effort to mitigate presented threats is to make all our data available for analysis on our website [27]. This includes the source code that categorizes test cases, the GUI models employed together with test case information.

V. CONCLUSIONS AND FUTURE WORK

In this paper we presented an initial study on what can be expected in long term GUI test case replayability in the case of complex open source software. We performed a cross-sectional evaluation where generated test cases were replayed on the immediately following studied version of the target application. We used this information to detail the results of the longitudinal evaluation where we performed a simulation of long term replayability for GUI test cases. We believe our results validate previous work in GUI test case maintenance [17]–[20] and we hope to fuel further work in the field.

Due to our promising initial results, our next goal is to conduct a comprehensive follow-up study employing a larger selection of target applications on other platforms such as .NET and SWT. This entails extending our software repository with new application versions and obtaining associated information regarding functionally equivalent GUI widgets. In addition, we aim to switch from *simulating* test case execution to running them by employing GUITAR's test runner component, thus eliminating one of the threats to the validity of the presented research.

¹⁵From version 0.6.7 it also has an *Options* window that could not be correctly ripped using GUIRipper.

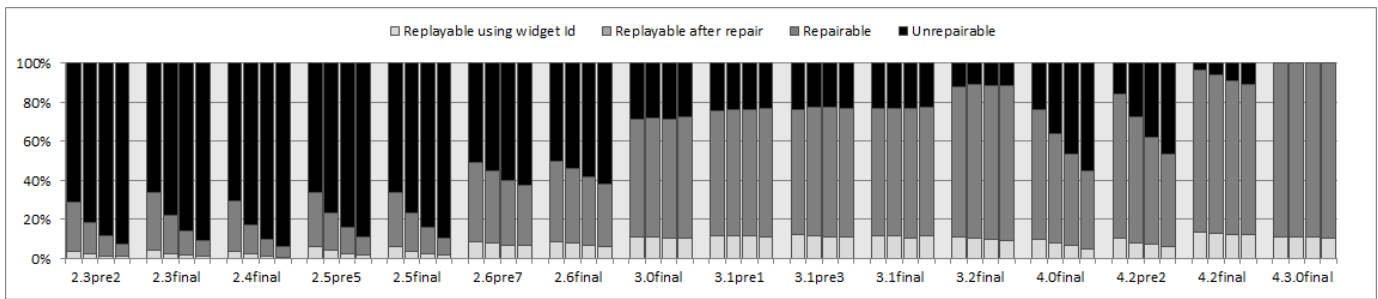


Fig. 4. Longitudinal jEdit test case replayability

A more distant avenue of research regards a comprehensive evaluation targeting event-driven systems beyond the desktop paradigm by including web and mobile applications. Theoretical advances in unified modelling of event driven software [28] together with GUITAR components targeting these platforms enable such a complex undertaking. We believe such an effort can lead to better understanding of GUI-driven software in a platform independent manner and enable the creation of unified testable models for such systems.

ACKNOWLEDGMENT

The author was supported by programs co-financed by The Sectoral Operational Programme Human Resources Development, Contract POS DRU 6/1.5/S/3 - "Doctoral studies: through science towards society"

REFERENCES

- [1] A. M. Memon, "A comprehensive framework for testing graphical user interfaces," Ph.D. dissertation, 2001, aAI3026063.
- [2] A. P. Nikolai, A. C. R. Paiva, N. Tillmann, J. C. P. Faria, and R. F. A. M., "Modeling and testing hierarchical guis," in *Proc.ASM05. Universit de Paris 12*, 2005, pp. 8–11.
- [3] A. C. R. P. Pimenta, "Automated specification-based testing of graphical user interfaces," Ph.D. dissertation, 2006.
- [4] A. Memon, "Gui ripping: Reverse engineering of graphical user interfaces for testing," in *In Proceedings of The 10th Working Conference on Reverse Engineering*, 2003, pp. 260–269.
- [5] Website., <http://guitar.sourceforge.net/>. Home of the GUITAR toolset.
- [6] A. Memon, A. Nagarajan, and Q. Xie, "Automating regression testing for evolving gui software," *Journal of Software Maintenance*, vol. 17, pp. 27–64, January 2005. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1062996.1062999>
- [7] D. Hackner and A. M. Memon, "Test case generator for GUITAR," in *ICSE '08: Research Demonstration Track: International Conference on Software Engineering*. Washington, DC, USA: IEEE Computer Society, 2008.
- [8] Q. Xie, "Developing cost-effective model-based techniques for gui testing," Ph.D. dissertation, College Park, MD, USA, 2006, aAI3241432.
- [9] Q. Xie and A. M. Memon, "Model-based testing of community-driven open-source gui applications," in *Proceedings of the 22nd IEEE International Conference on Software Maintenance*. Washington, DC, USA: IEEE Computer Society, 2006, pp. 145–154. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1172962.1172990>
- [10] M. L. Hammontree, J. J. Hendrickson, and B. W. Hensley, "Integrated data capture and analysis tools for research and testing on graphical user interfaces," in *Proceedings of the SIGCHI conference on Human factors in computing systems*, ser. CHI '92. New York, NY, USA: ACM, 1992, pp. 431–432.
- [11] A. Memon and et al., "What test oracle should i use for effective gui testing?" in *PROC. IEEE INTERNATIONAL CONFERENCE ON AUTOMATED SOFTWARE ENGINEERING (ASE'03)*. IEEE Computer Society Press, 2003, pp. 164–173.
- [12] J. Takahashi, "An automated oracle for verifying gui objects," *SIGSOFT Softw. Eng. Notes*, vol. 26, pp. 83–88, July 2001. [Online]. Available: <http://doi.acm.org/10.1145/505482.505494>
- [13] T. Ostrand, A. Anodide, H. Foster, and T. Goradia, "A visual test development environment for gui systems," *SIGSOFT Softw. Eng. Notes*, vol. 23, pp. 82–92, March 1998. [Online]. Available: <http://doi.acm.org/10.1145/271775.271793>
- [14] H. Hemmati, A. Arcuri, and L. Briand, "Reducing the cost of model-based testing through test case diversity," in *Proceedings of the 22nd IFIP WG 6.1 international conference on Testing software and systems*, ser. ICTSS'10. Berlin, Heidelberg: Springer-Verlag, 2010, pp. 63–78.
- [15] J. Jacky, M. Veanes, C. Campbell, and W. Schulte, *Model-Based Software Testing and Analysis with C#*, 1st ed.
- [16] J. L. Silva, J. C. Campos, and A. C. R. Paiva, "Model-based user interface testing with spec explorer and concurtasktrees," *Electron. Notes Theor. Comput. Sci.*, vol. 208, pp. 77–93, April 2008.
- [17] A. M. Memon, "Automatically repairing event sequence-based gui test suites for regression testing," *ACM Trans. Softw. Eng. Methodol.*, vol. 18, pp. 4:1–4:36, November 2008. [Online]. Available: <http://doi.acm.org/10.1145/1416563.1416564>
- [18] S. Huang, M. B. Cohen, and A. M. Memon, "Repairing gui test suites using a genetic algorithm," in *Proceedings of the 2010 Third International Conference on Software Testing, Verification and Validation*, ser. ICST '10. Washington, DC, USA: IEEE Computer Society, 2010, pp. 245–254. [Online]. Available: <http://dx.doi.org/10.1109/ICST.2010.39>
- [19] S. McMaster and A. M. Memon, "An extensible heuristic-based framework for gui test case maintenance," in *Proceedings of the IEEE International Conference on Software Testing, Verification, and Validation Workshops*. Washington, DC, USA: IEEE Computer Society, 2009, pp. 251–254.
- [20] A.-J. Molnar, "A heuristic process for GUI widget matching across application versions," *Annales Universitatis. Scientiarum Budapestensis, Sectio Computatorica*, vol. 36, pp. 255–275, 2012.
- [21] A. M. Memon, "An event-flow model of gui-based applications for testing," *Software Testing, Verification and Reliability*, vol. 17, no. 3, pp. 137–157, 2007.
- [22] P. A. Brooks and A. M. Memon, "Automated gui testing guided by usage profiles," in *Proceedings of the twenty-second IEEE/ACM international conference on Automated software engineering*, ser. ASE '07. New York, NY, USA: ACM, 2007, pp. 333–342.
- [23] Q. Xie and A. M. Memon, "Studying the characteristics of a 'good' GUI test suite," in *Proceedings of the 17th IEEE International Symposium on Software Reliability Engineering (ISSRE 2006)*. IEEE Computer Society Press, Nov. 2006.
- [24] —, "Using a pilot study to derive a GUI model for automated testing," *ACM Trans. Softw. Eng. Methodol.*, vol. 18, November 2008.
- [25] A. C. R. Paiva, J. C. P. Faria, and R. F. A. M. Vidal, "Towards the integration of visual and formal models for gui testing," *Electr. Notes Theor. Comput. Sci.*, vol. 190, no. 2, pp. 99–111, 2007.
- [26] A.-J. Molnar, "A software repository and toolset for empirical research," *Studia Informatica UBB*, vol. LVII, no. 1, pp. 73–88, March 2012.
- [27] Website., <https://sourceforge.net/projects/javaset> (Home of our GUI toolset).
- [28] R. Bryce, S. Sampath, and A. M. Memon, "Developing a single model and test prioritization strategies for event-driven software," *IEEE Transactions on Software Engineering*, 2011.